

# CS391L: Machine Learning WB

Language Modeling: RNNs, Attention and Transformers

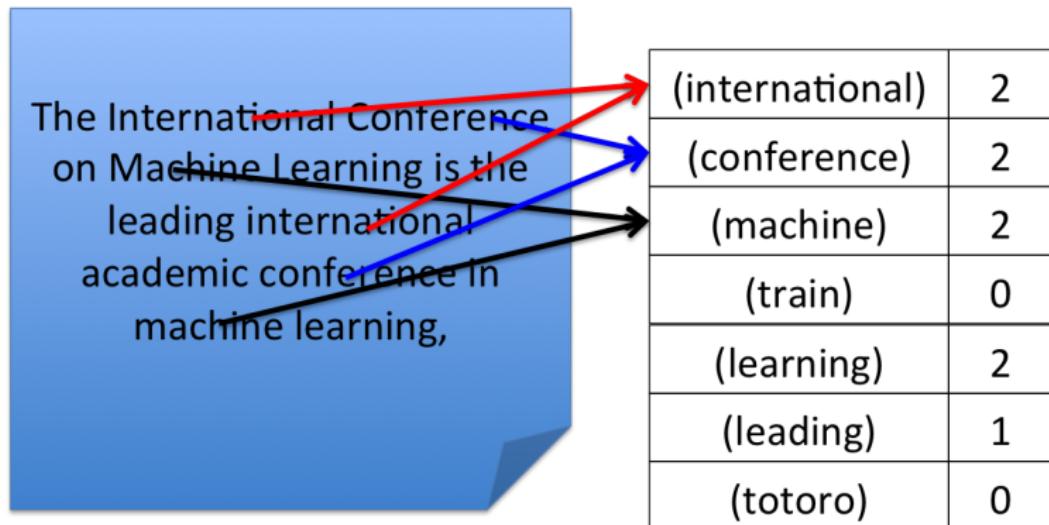
Inderjit Dhillon  
UT Austin

April 1, 2025

# Representation for sentence/document

# Bag of Word

- A classical way to represent NLP data
- Each sentence (or document) is represented by a  $d$ -dimensional vector  $\mathbf{x}$ , where  $x_i$  is number of occurrences of word  $i$ .



number of features = number of potential words (very large)

# Feature generation for documents

- Bag of  $n$ -gram features ( $n = 2$ ):

The International Conference on Machine Learning is the leading international academic conference in machine learning,

(international)	2
(conference)	2
(machine)	2
(train)	0
(learning)	2
(leading)	1
(totoro)	0

(international conference)	1
(machine learning)	2
(leading international)	1
(totoro tiger)	0
(tiger woods)	0
(international academic)	1
(international academic)	1

10,000 words  $\Rightarrow 10,000^2$  potential features

# Data Matrix (document)

- Use the bag-of-words matrix or the normalized version (TF-IDF) for a dataset (denoted by  $D$ ):

$$\text{tfidf}(\text{doc}, \text{word}, D) = \text{tf}(\text{doc}, \text{word}) \cdot \text{idf}(\text{word}, D)$$

- $\text{tf}(\text{doc}, \text{word})$ : term frequency

(word count in the document)/(total number of terms in the document)

- $\text{idf}(\text{word}, \text{Dataset})$ : inverse document frequency

$1 + \log((\text{Number of documents}) / (\text{Number of documents with this word}))$

	angeles	los	new	post	times	york
d1	0	0	1	0	1	1
d2	0	0	1	1	0	1
d3	1	1	0	0	1	0

tf-idf

	angeles	los	new	post	times	york
d1	0	0	0.584	0	0.584	0.584
d2	0	0	0.584	1.584	0	0.584
d3	1.584	1.584	0	0	0.584	0

## Bag of words + linear model

- Example: text classification (e.g., sentiment prediction, review score prediction)
- Linear model:  $y \approx \text{sign}(\mathbf{w}^T \mathbf{x})$   
(e.g., by linear SVM/logistic regression)
- $w_i$ : the “contribution” of each word

# Bag of words + Fully connected network

- $f(\mathbf{x}) = W_L \sigma(W_{L-1} \cdots \sigma(W_0 \mathbf{x}))$
- The first layer  $W_0$  is a  $d_1 \times d$  matrix:
  - Each column  $\mathbf{w}_i$  is a  $d_1$ -dimensional representation of  $i$ -th word ([word embedding](#))
  - $W_0 \mathbf{x} = x_1 \mathbf{w}_1 + x_2 \mathbf{w}_2 + \cdots + x_d \mathbf{w}_d$  is a [linear combination](#) of these vectors
  - $W_0$  is also called the [word embedding matrix](#)
  - Final prediction can be viewed as an  $L - 1$  layer network on  $W_0 \mathbf{x}$  (average of word embeddings)
- Not capturing the sequence information in language

# Recurrent Neural Network

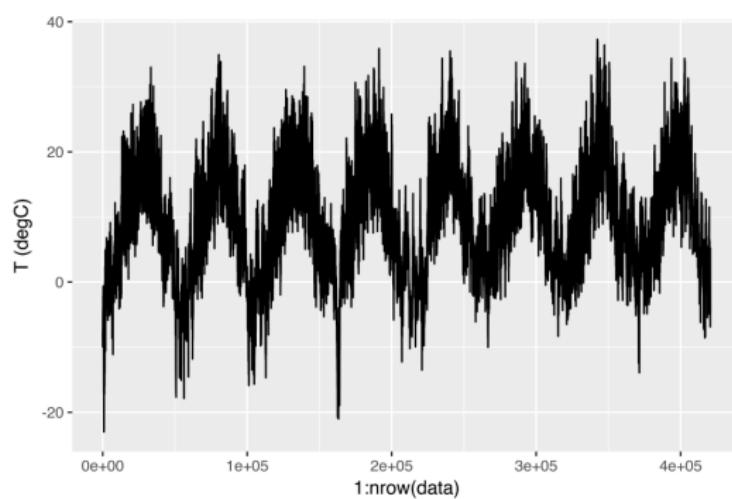
# Time Series/Sequence Data

- Input:  $\{x_1, x_2, \dots, x_T\}$ 
  - Each  $x_t$  is the feature at time step  $t$
  - Each  $x_t$  can be an  $d$ -dimensional vector
- Output:  $\{y_1, y_2, \dots, y_T\}$ 
  - Each  $y_t$  is the output at step  $t$
  - Multi-class output or Regression output:

$$y_t \in \{1, 2, \dots, L\} \text{ or } y_t \in \mathbb{R}$$

# Example: Time Series Prediction

- Climate Data:
  - $x_t$ : temperature at time  $t$
  - $y_t$ : temperature (or temperature change) at time  $t + 1$
- Stock Price: Predicting stock price



# Example: Language Modeling

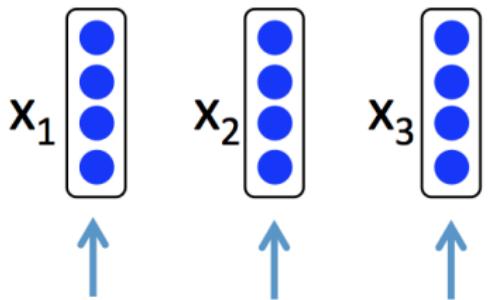
The cat is ?

## Example: Language Modeling

The cat is ?

- $x_t$ : one-hot encoding to represent the word at step  $t$   
([0, ..., 0, 1, 0, ..., 0])
- $y_t$ : the next word

$$y_t \in \{1, \dots, V\} \quad V: \text{Vocabulary size}$$

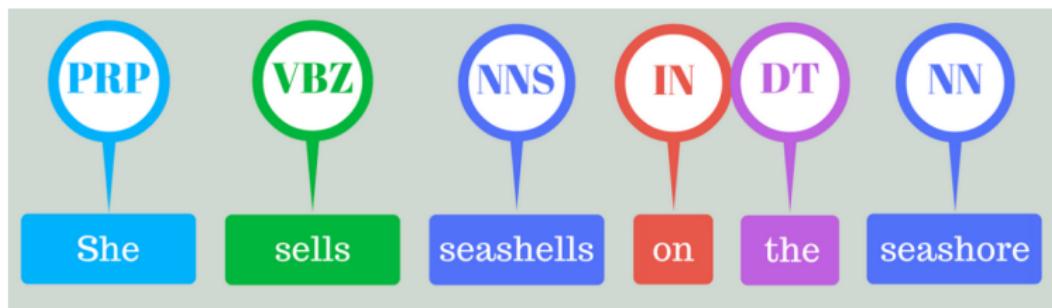


The cat is ?

$y_1$        $y_2$        $y_3$

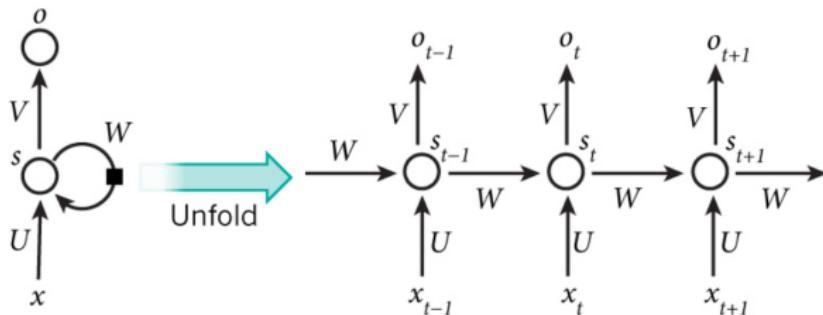
# Example: POS Tagging

- Part of Speech Tagging:  
Labeling words with their Part-Of-Speech (Noun, Verb, Adjective, ...)
- $x_t$ : a vector to represent the word at step  $t$
- $y_t$ : label of word  $t$



picture from <https://medium.com/analytics-vidhya/pos-tagging-using-conditional-random-fields-92077e5eaa31>

# Recurrent Neural Network (RNN)



- $x_t$ :  $t$ -th input
- $s_t$ : hidden state at time  $t$  ("memory" of the network)

$$s_t = f(Ux_t + Ws_{t-1})$$

$W$ : transition matrix,  $U$  : word embedding matrix

$s_0$  usually set to be 0

- Predicted output at time  $t$ :

$$o_t = \arg \max_i (Vs_t)_i$$

# Recurrent Neural Network (RNN)

- Training: Find  $U, W, V$  to minimize empirical loss:
- Loss of a sequence:

$$\sum_{t=1}^T \text{loss}(V\mathbf{s}_t, y_t)$$

( $\mathbf{s}_t$  is a function of  $U, W, V$ )

- Loss on the whole dataset:  
Average loss over all sequences
- Solved by SGD/Adam

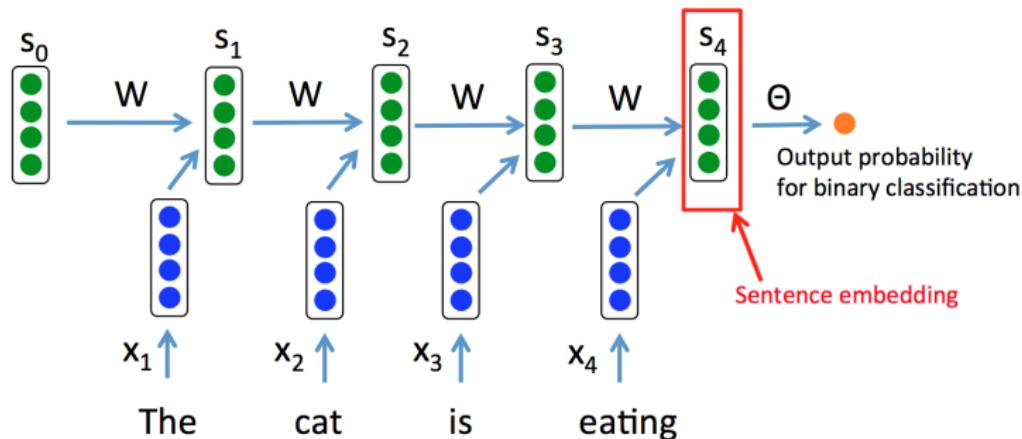
# RNN: Text Classification

- Not necessary to output at each step
- Text Classification:

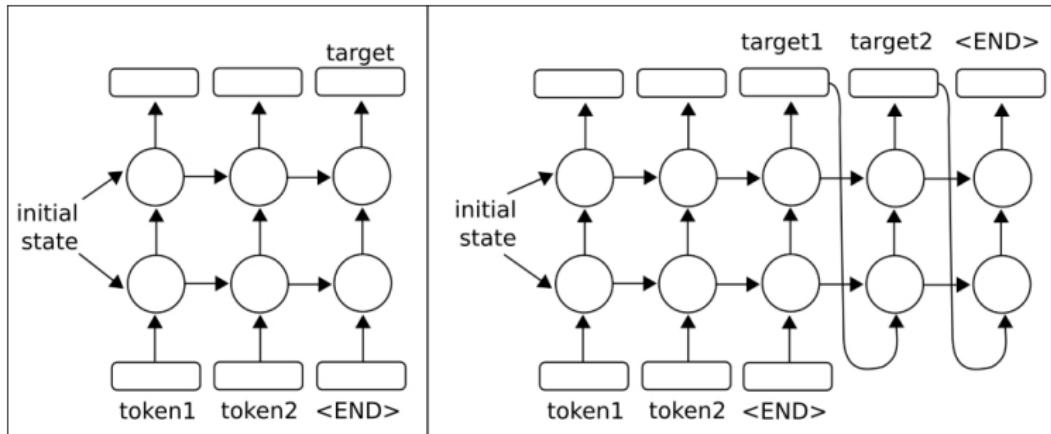
sentence → category

Output only at the final step

- Model: add a fully connected network to the final embedding



# Multi-layer RNN



(Figure from

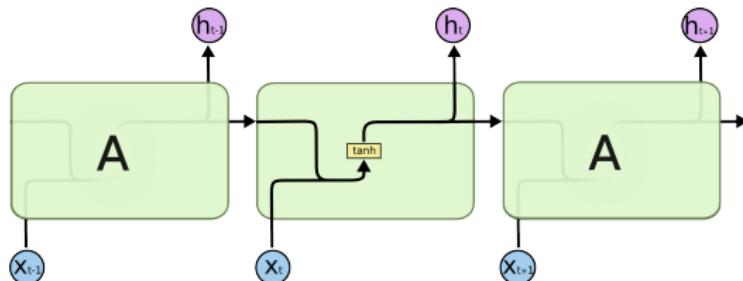
[https://subscription.packtpub.com/book/big\\_data\\_and\\_business\\_intelligence](https://subscription.packtpub.com/book/big_data_and_business_intelligence))

# Problems with Classical RNN

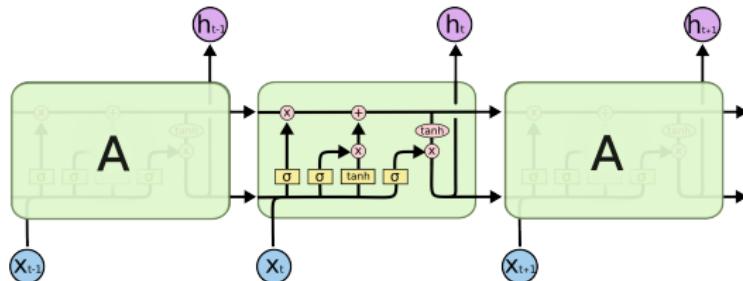
- Hard to capture long-term dependencies
- Hard to solve (vanishing gradient problem)
- Solution:
  - LSTM (Long Short Term Memory networks)
  - GRU (Gated Recurrent Unit)
  - ...

# LSTM

- RNN:



- LSTM:



Neural Network Layer

Pointwise Operation

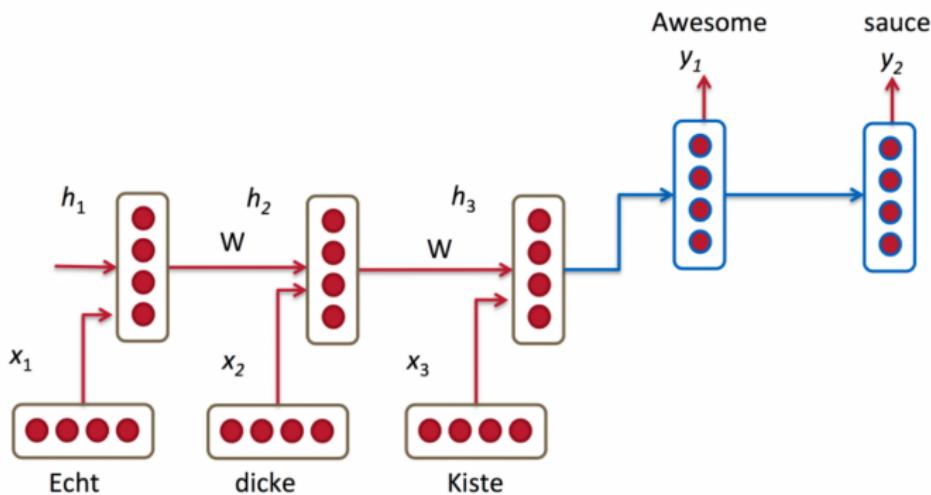
Vector Transfer

Concatenate

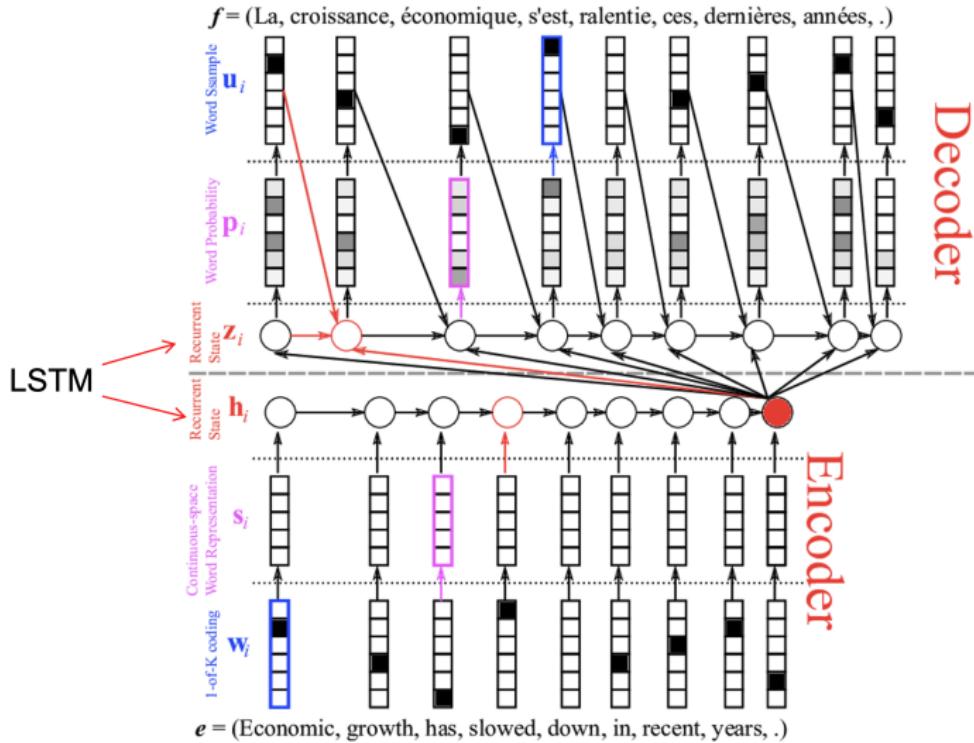
Copy

# Neural Machine Translation (NMT)

- Output the translated sentence given an input sentence
- Training data: a set of input-output pairs (supervised setting)
- Encoder-decoder approach:
  - Encoder: Use (RNN/LSTM) to encode the input sentence into a latent vector
  - Decoder: Use (RNN/LSTM) to generate a sentence based on the latent vector



# Neural Machine Translation



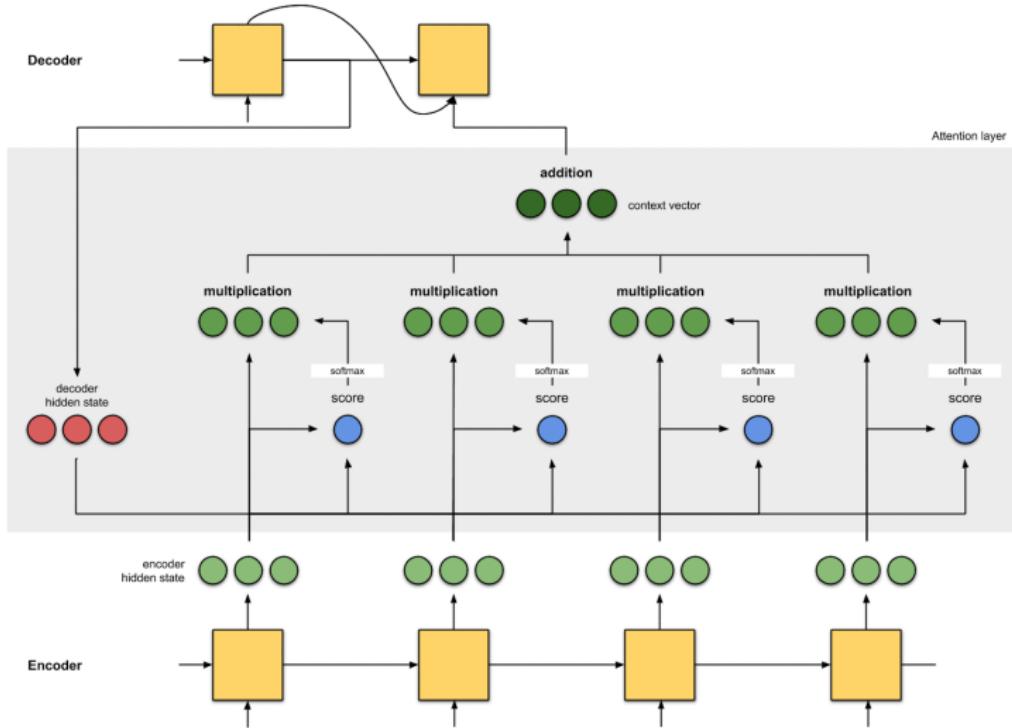
# Attention in NMT

- Usually, each output word is only related to a subset of input words (e.g., for machine translation)
- Let  $\mathbf{u}$  be the current decoder latent state  $\mathbf{v}_1, \dots, \mathbf{v}_n$  be the latent state for each input word
- Compute the weight of each state by

$$\mathbf{p} = \text{Softmax}(\mathbf{u}^T \mathbf{v}_1, \dots, \mathbf{u}^T \mathbf{v}_n)$$

- Compute the context vector by  $V\mathbf{p} = p_1 \mathbf{v}_1 + \dots + p_n \mathbf{v}_n$

# Attention in NMT

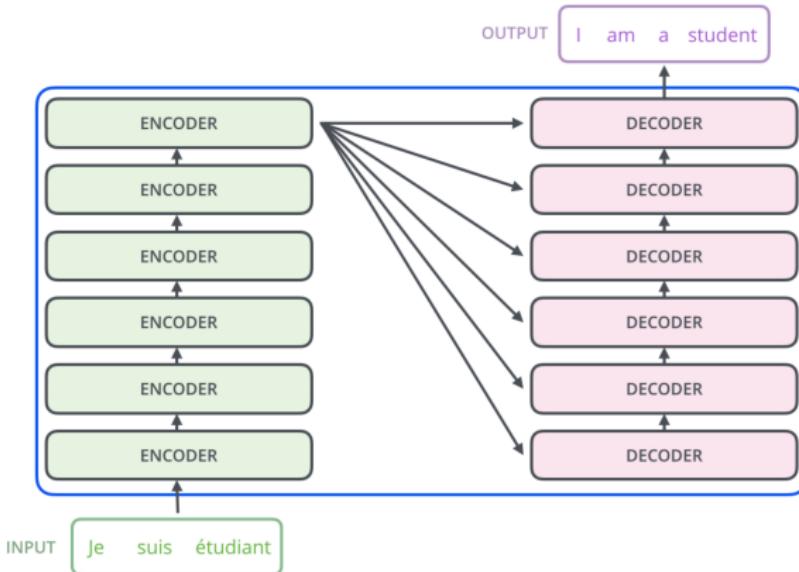


(Figure from <https://towardsdatascience.com/>)

# Transformer

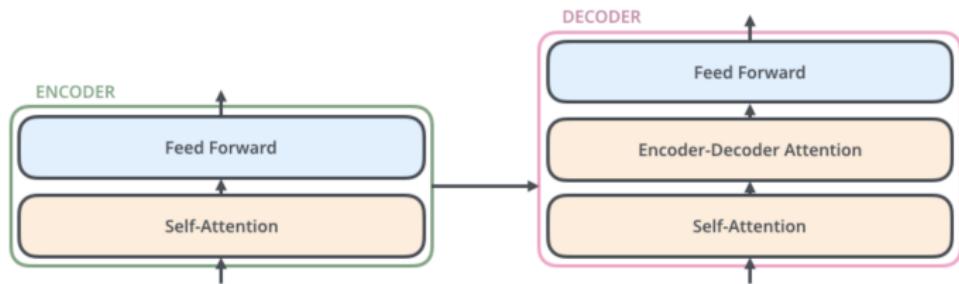
# Transformer

- An architecture that replies entirely on attention without using CNN/RNN
- Proposed in “Attention Is All You Need” (Vaswani et al., 2017)
- Initially used for neural machine translation



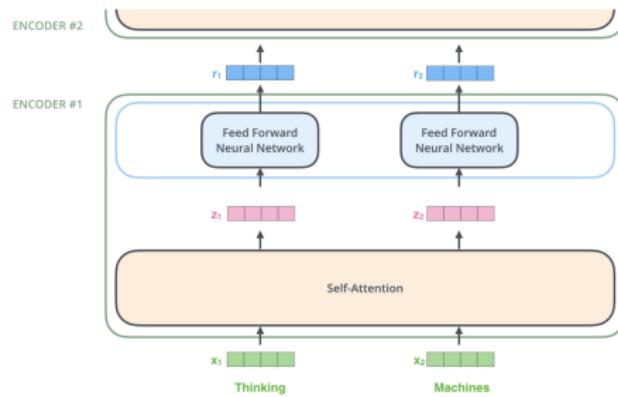
# Encoder and Decoder

- Self attention layer: the main architecture used in Transformer
- Decoder: will have another attention layer to help it focus on relevant parts of the input sentence.



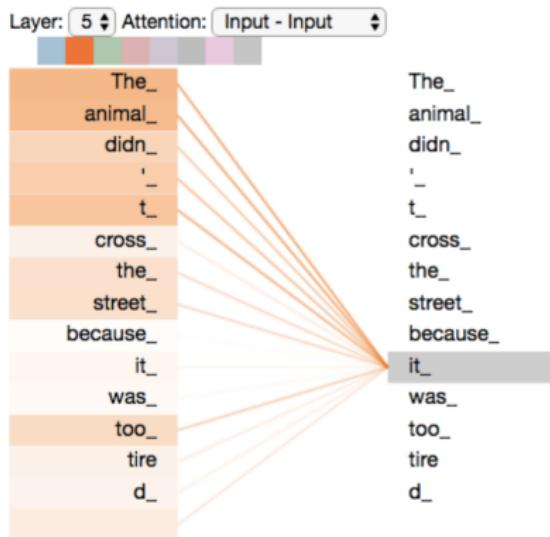
# Encoder

- Each word has a corresponding “latent vector” (initially the word embedding for each word)
- Each layer of encoder:
  - Receive a list of vectors as input
  - Passing these vectors to a **self-attention** layer
  - Then passing them into a feed-forward layer
  - Output a list of vectors



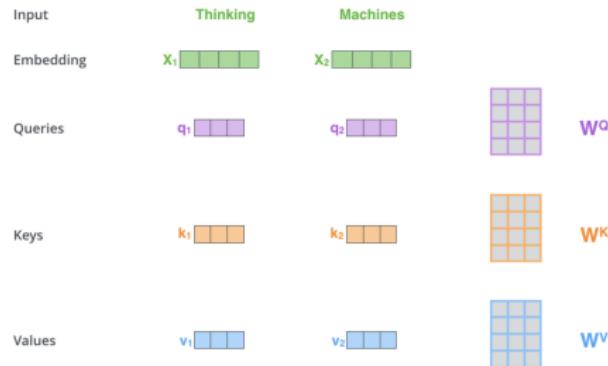
# Self-attention layer

- Main idea: The actual meaning of each word may be related to other words in the sentence
- The actual meaning (latent vector) of each word is a weighted (attention) combination of other words (latent vectors) in the sentences



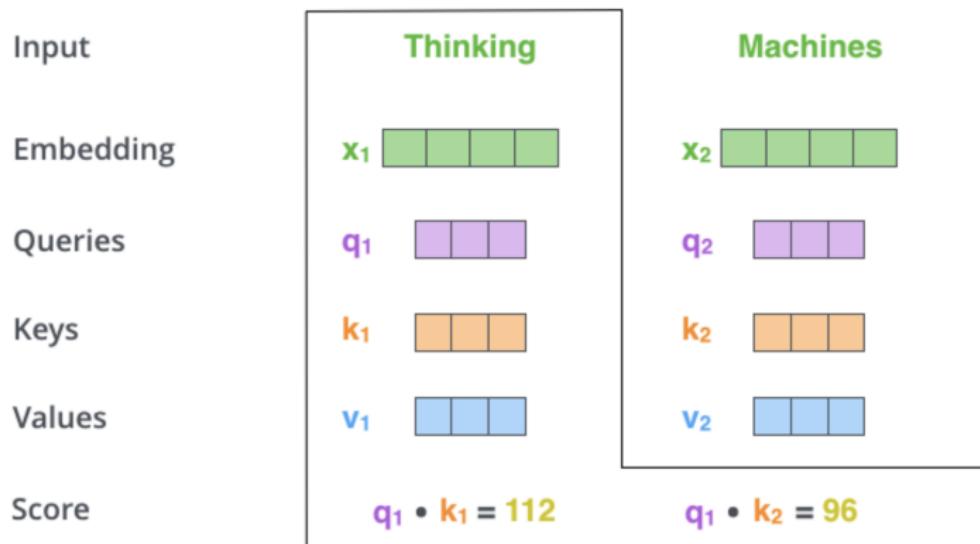
# Self-attention layer

- Input latent vectors:  $x_1, \dots, x_n$
- Self-attention parameters:  $W^Q, W^K, W^V$  (weights for query, key, value)
- For each word  $i$ , compute
  - Query vector:  $q_i = x_i W^Q$
  - Key vector:  $k_i = x_i W^K$
  - Value vector:  $v_i = x_i W^V$



# Self-attention layer

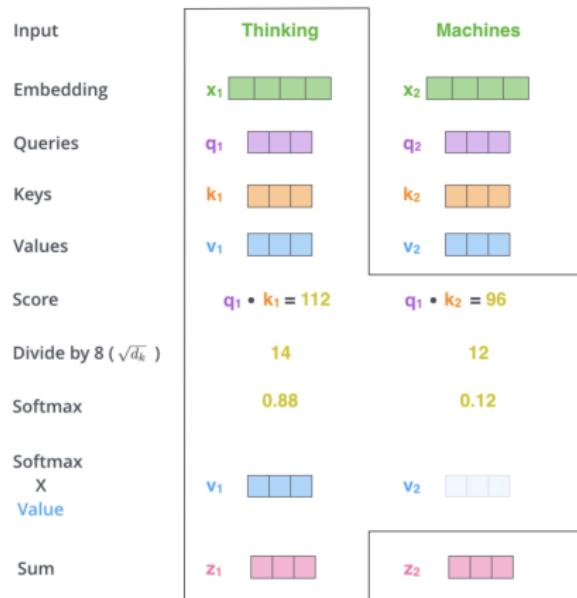
- For each word  $i$ , compute the scores to determine how much focus to place on other input words
  - The **attention** score for word  $j$  to word  $i$ :  $q_i^T k_j$



# Self-attention layer

- For each word  $i$ , the output vector

$$\sum_j s_{ij} \mathbf{v}_j, \quad \mathbf{s}_i = \text{softmax}(\mathbf{q}_i^T \mathbf{k}_1, \dots, \mathbf{q}_i^T \mathbf{k}_n)$$



# Matrix form

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V, \quad Z = \text{softmax}(QK^T)V$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} W^Q \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} Q \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

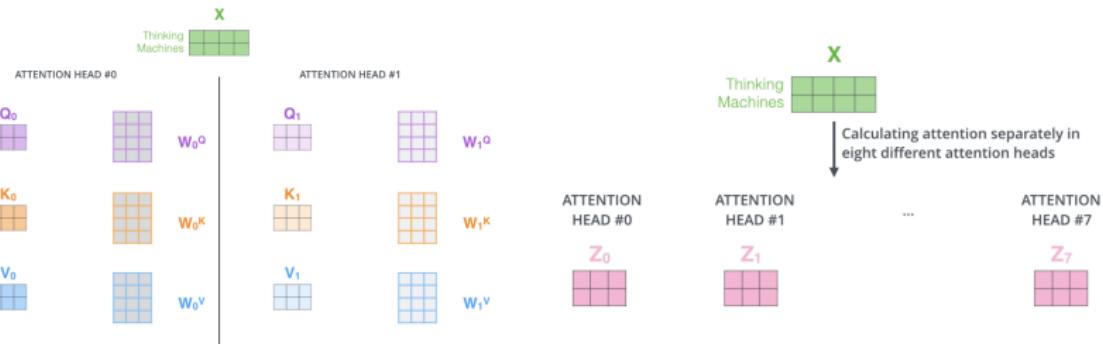
$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} W^K \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} K \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} W^V \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} V \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

$$\text{softmax}\left(\frac{\begin{matrix} Q \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix} \times \begin{matrix} K^T \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) = \begin{matrix} Z \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

# Multiple heads

- Multi-headed attention: use multiple set of (*key*, *value*, *query*) weights
- Each head will output a vector  $Z_i$



# Multiply with weight matrix to reshape

- Gather all the outputs  $Z_1, \dots, Z_k$
- Multiply with a weight matrix to reshape
- Then pass to the next fully connected layer

1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^o$  that was trained jointly with the model

$X$

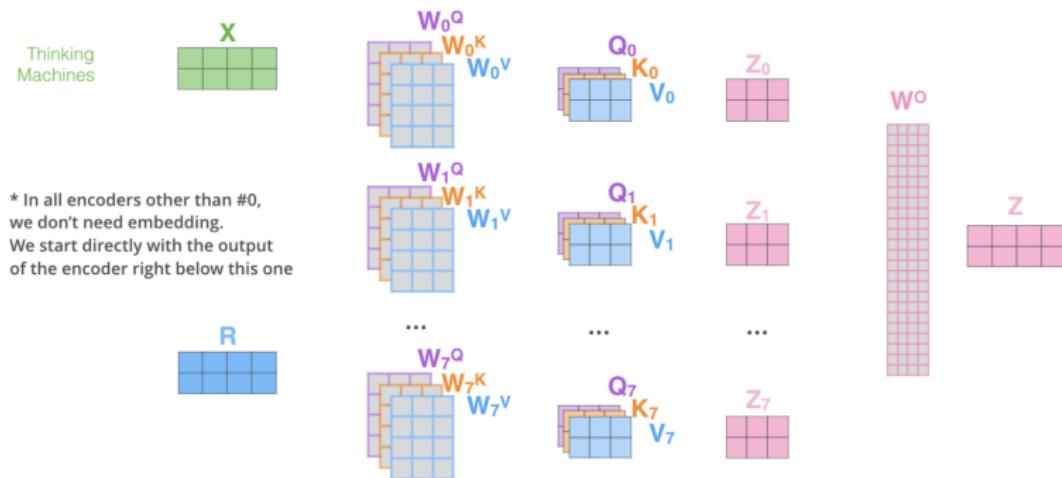


3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \begin{matrix} \text{---} & \text{---} & \text{---} & \text{---} \end{matrix} \end{matrix}$$

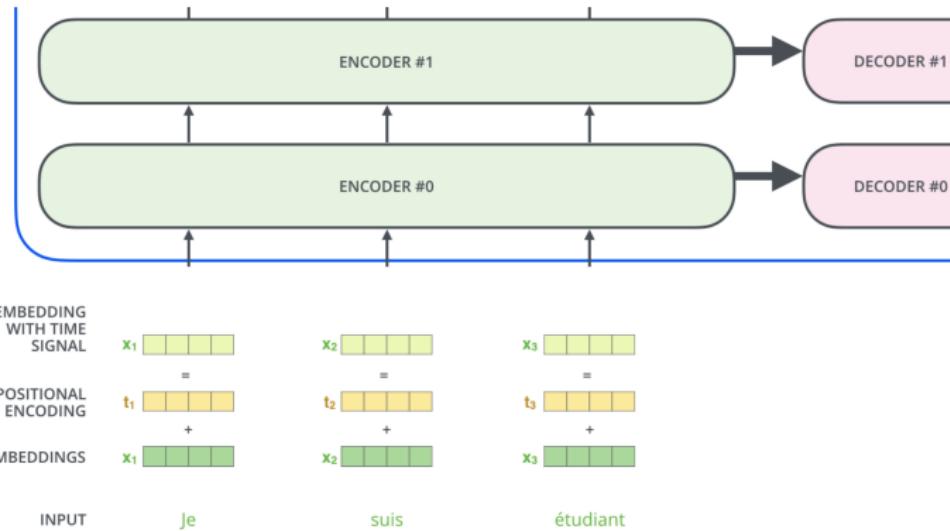
# Overall architecture

- 1) This is our input sentence\*  
Thinking Machines
- 2) We embed each word\*
- 3) Split into 8 heads.  
We multiply  $X$  or  $R$  with weight matrices
- 4) Calculate attention using the resulting  $Q/K/V$  matrices
- 5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^o$  to produce the output of the layer



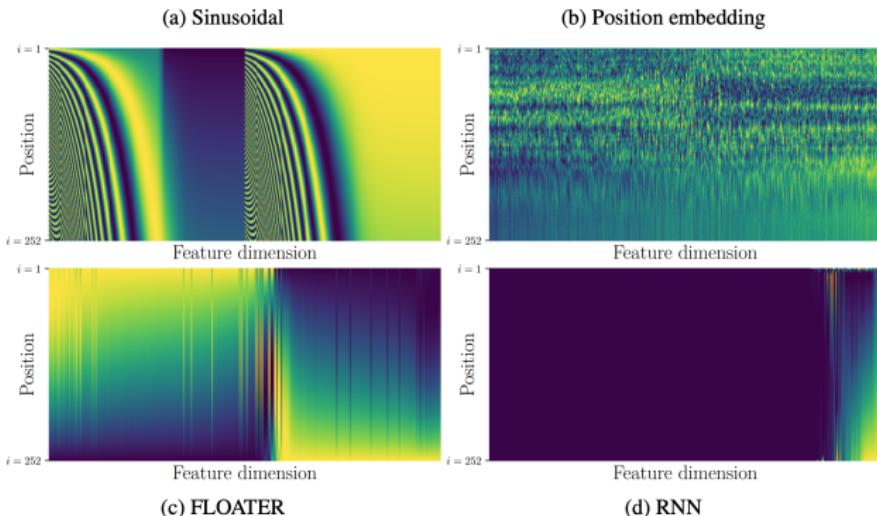
# Position encoding

- The above architecture ignores the sequential information
- Add a position encoding vector to each  $x_i$  (according to  $i$ )



# Types of position encoding

- Sin/cosine functions with different wavelengths (used in the original Transformer)
  - smooth, parameter-free, inductive
- Position embedding: learn a latent vector for each position
  - non-smooth, data-driven (learnable), non-inductive
- Neural ODE embedding (Liu et al., 2020)
  - smooth, data-driven (learnable), inductive



# Conclusions

- A brief introduction of RNN and Transformer.

Questions?